

# Practical Issues

or avoiding machine learning going mad

JJ Valletta

March 6, 2015

[www.exeter.ac.uk/as/rdp/](http://www.exeter.ac.uk/as/rdp/)

# Overview

- Machine learning gone mad
- Laws of data analysis
- A few important tips
- Which machine learning algorithm should I use?
- Validating models

# Overview

- Machine learning gone mad
- Laws of data analysis
- A few important tips
- Which machine learning algorithm should I use?
- Validating models

# Overview

- Machine learning gone mad
- Laws of data analysis
- A few important tips
- Which machine learning algorithm should I use?
- Validating models

# Overview

- Machine learning gone mad
- Laws of data analysis
- A few important tips
- Which machine learning algorithm should I use?
- Validating models

# Overview

- Machine learning gone mad
- Laws of data analysis
- A few important tips
- Which machine learning algorithm should I use?
- Validating models

# Overview

- Machine learning gone mad
- Laws of data analysis
- A few important tips
- Which machine learning algorithm should I use?
- Validating models

# Machine learning gone mad - The US tank experiment

- 1980s: Pentagon got excited about artificial neural networks
  - Wanted a classifier to detect whether a tank is hiding behind trees
  - Collected photos of trees with/without tanks hiding in them
  - Trained neural network performed *excellently* on the testing dataset  
**Champagne to all the scientists!**





# Machine learning gone mad - The US tank experiment

- 1980s: Pentagon got excited about artificial neural networks
- Wanted a classifier to detect whether a tank is hiding behind trees
- Collected photos of trees with/without tanks hiding in them
- Trained neural network performed *excellently* on the testing dataset  
**Champagne to all the scientists!**



# Machine learning gone mad - The US tank experiment

- 1980s: Pentagon got excited about artificial neural networks
- Wanted a classifier to detect whether a tank is hiding behind trees
- Collected photos of trees with/without tanks hiding in them
- Trained neural network performed *excellently* on the testing dataset  
**Champagne to all the scientists!**



# Machine learning gone mad - The US tank experiment

- 1980s: Pentagon got excited about artificial neural networks
- Wanted a classifier to detect whether a tank is hiding behind trees
- Collected photos of trees with/without tanks hiding in them
- Trained neural network performed *excellently* on the testing dataset  
**Champagne to all the scientists!**



# Machine learning gone mad - The US tank experiment

- Another set of tank/no tank photos was commissioned
- The classifier was now useless and no better than randomly guessing
- Someone then noted the following on the original dataset:  
no tank: *all* photos taken on a sunny, **blue** skies day  
tank: *all* photos taken on a cloudy, **grey** skies day
- Neural network had learnt whether it was a sunny or cloudy day  
**God bless the United States of America!**



# Machine learning gone mad - The US tank experiment

- Another set of tank/no tank photos was commissioned
- The classifier was now useless and no better than randomly guessing
- Someone then noted the following on the original dataset:
  - no tank: *all* photos taken on a sunny, **blue** skies day
  - tank: *all* photos taken on a cloudy, **grey** skies day
- Neural network had learnt whether it was a sunny or cloudy day  
**God bless the United States of America!**



# Machine learning gone mad - The US tank experiment

- Another set of tank/no tank photos was commissioned
- The classifier was now useless and no better than randomly guessing
- Someone then noted the following on the original dataset:
  - no tank:** *all* photos taken on a sunny, **blue** skies day
  - tank:** *all* photos taken on a cloudy, **grey** skies day
- Neural network had learnt whether it was a sunny or cloudy day  
**God bless the United States of America!**



# Machine learning gone mad - The US tank experiment

- Another set of tank/no tank photos was commissioned
- The classifier was now useless and no better than randomly guessing
- Someone then noted the following on the original dataset:
  - no tank:** *all* photos taken on a sunny, **blue** skies day
  - tank:** *all* photos taken on a cloudy, **grey** skies day
- Neural network had learnt whether it was a sunny or cloudy day  
**God bless the United States of America!**



# Machine learning gone mad - Google flu trends

- 2008: Google decided to showcase the power of big data
  - Wanted a predictive model to estimate influenza activity
  - Training data were queries containing terms such as cough and fever
  - Used IP address to break data by states
  - Outcome (label) was influenza-like illness (ILI) physician visits collected by CDC (Centers for Disease Control and Prevention)
  - Model was successful, it predicted a spike in the mid-Atlantic region of the US two weeks prior to the CDC
- Google is great!**



# Machine learning gone mad - Google flu trends

- 2008: Google decided to showcase the power of big data
  - Wanted a predictive model to estimate influenza activity
  - Training data were queries containing terms such as cough and fever
  - Used IP address to break data by states
  - Outcome (label) was influenza-like illness (ILI) physician visits collected by CDC (Centers for Disease Control and Prevention)
  - Model was successful, it predicted a spike in the mid-Atlantic region of the US two weeks prior to the CDC
- Google is great!**

# Machine learning gone mad - Google flu trends

- 2008: Google decided to showcase the power of big data
  - Wanted a predictive model to estimate influenza activity
  - Training data were queries containing terms such as cough and fever
  - Used IP address to break data by states
  - Outcome (label) was influenza-like illness (ILI) physician visits collected by CDC (Centers for Disease Control and Prevention)
  - Model was successful, it predicted a spike in the mid-Atlantic region of the US two weeks prior to the CDC
- Google is great!**

# Machine learning gone mad - Google flu trends

- 2008: Google decided to showcase the power of big data
  - Wanted a predictive model to estimate influenza activity
  - Training data were queries containing terms such as cough and fever
  - Used IP address to break data by states
  - Outcome (label) was influenza-like illness (ILI) physician visits collected by CDC (Centers for Disease Control and Prevention)
  - Model was successful, it predicted a spike in the mid-Atlantic region of the US two weeks prior to the CDC
- Google is great!**

# Machine learning gone mad - Google flu trends

- 2008: Google decided to showcase the power of big data
  - Wanted a predictive model to estimate influenza activity
  - Training data were queries containing terms such as cough and fever
  - Used IP address to break data by states
  - Outcome (label) was influenza-like illness (ILI) physician visits collected by CDC (Centers for Disease Control and Prevention)
  - Model was successful, it predicted a spike in the mid-Atlantic region of the US two weeks prior to the CDC
- Google is great!

# Machine learning gone mad - Google flu trends

- 2008: Google decided to showcase the power of big data
- Wanted a predictive model to estimate influenza activity
- Training data were queries containing terms such as cough and fever
- Used IP address to break data by states
- Outcome (label) was influenza-like illness (ILI) physician visits collected by CDC (Centers for Disease Control and Prevention)
- Model was successful, it predicted a spike in the mid-Atlantic region of the US two weeks prior to the CDC

**Google is great!**

# Machine learning gone mad - Google flu trends

google.org Flu Trends

[Google.org home](#)

[Dengue Trends](#)

**Flu Trends**

[Home](#)

United States

National

[Download data](#)

[How does this work?](#)

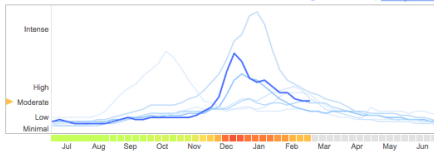
[FAQ](#)

## Explore flu trends - United States

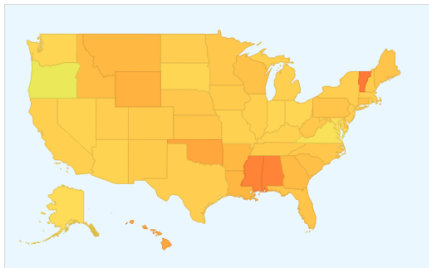
We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. [Learn more »](#)

### National

● 2014-2015 ● Past years ▼



States | [Cities](#) (Experimental)



Estimates were made using a model that proved accurate when compared to historic official flu activity data. Data current through March 3, 2015.

### BIG DATA

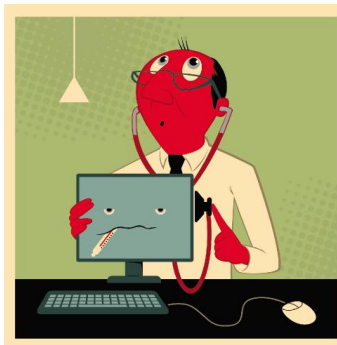
## The Parable of Google Flu: Traps in Big Data Analysis

David Lazer,<sup>1,2\*</sup> Ryan Kennedy,<sup>1,3,4</sup> Gary King,<sup>2</sup> Alessandro Vespignani<sup>3,5,6</sup>

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (1, 2). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data (3, 4), what lessons can we draw from this error?

The problems we identify are not limited to GFT. Research on whether search or social media can predict  $x$  has become common-

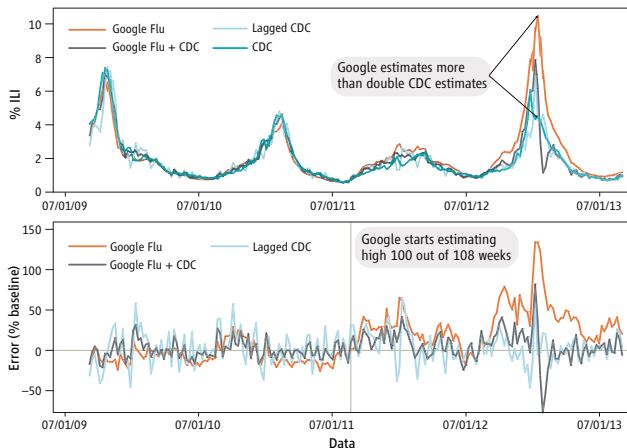
Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.



the algorithm in 2009, and this model has run ever since, with a few changes announced in October 2013 (10, 15).

Although not widely reported until 2013, the new GFT has been persistently overestimating flu prevalence for a much longer time. GFT also missed by a very large margin in the 2011–2012 flu season and has missed high for 100 out of 108 weeks starting with August 2011 (see the graph). These errors are not randomly distributed. For example, last week's errors predict this week's errors (temporal autocorrelation), and the direction and magnitude of error varies with the time of year (seasonality). These patterns mean that GFT overlooks considerable information that could be extracted by traditional statistical methods.

# Machine learning gone mad - Google flu trends



**GFT overestimation.** GFT overestimated the prevalence of flu in the 2012–2013 season and overshot the actual level in 2011–2012 by more than 50%. From 21 August 2011 to 1 September 2013, GFT reported overly high flu prevalence 100 out of 108 weeks. **(Top)** Estimates of doctor visits for ILI. "Lagged CDC" incorporates 52-week seasonality variables with lagged CDC data. "Google Flu + CDC" combines GFT, lagged CDC estimates, lagged error of GFT estimates, and 52-week seasonality variables. **(Bottom)** Error [as a percentage  $\{[(\text{Non-CDC estimate}) - (\text{CDC estimate})] / (\text{CDC estimate})\}]$ . Both alternative models have much less error than GFT alone. Mean absolute error (MAE) during the out-of-sample period is 0.486 for GFT, 0.311 for lagged CDC, and 0.232 for combined GFT and CDC. All of these differences are statistically significant at  $P < 0.05$ . See SM.



# Google flu trends, so what went wrong?

- Google had to mine 50 million search terms that best correlated with 1,152 “true” data points collected by CDC (Centers for Disease Control and Prevention). They retained 45 queries
- Winter is coming: Correlation is not causation
- Are these correlates stable and comparable over time?
- Google's search algorithm changes very frequently. Autosuggest feature might have implicitly caused more people to search flu-related terms
- Is it time to recalibrate the model and/or hybridise both data sources?
- In 2014, Google retrained and updated the model significantly  
**Will wait and see!**

# Google flu trends, so what went wrong?

- Google had to mine 50 million search terms that best correlated with 1,152 “true” data points collected by CDC (Centers for Disease Control and Prevention). They retained 45 queries
- Winter is coming: Correlation is not causation
  - Are these correlates stable and comparable over time?
  - Google's search algorithm changes very frequently. Autosuggest feature might have implicitly caused more people to search flu-related terms
  - Is it time to recalibrate the model and/or hybridise both data sources?
  - In 2014, Google retrained and updated the model significantly  
**Will wait and see!**

# Google flu trends, so what went wrong?

- Google had to mine 50 million search terms that best correlated with 1,152 “true” data points collected by CDC (Centers for Disease Control and Prevention). They retained 45 queries
- Winter is coming: Correlation is not causation
- Are these correlates stable and comparable over time?
- Google's search algorithm changes very frequently. Autosuggest feature might have implicitly caused more people to search flu-related terms
- Is it time to recalibrate the model and/or hybridise both data sources?
- In 2014, Google retrained and updated the model significantly  
**Will wait and see!**

# Google flu trends, so what went wrong?

- Google had to mine 50 million search terms that best correlated with 1,152 “true” data points collected by CDC (Centers for Disease Control and Prevention). They retained 45 queries
- Winter is coming: Correlation is not causation
- Are these correlates stable and comparable over time?
- Google’s search algorithm changes very frequently. Autosuggest feature might have implicitly caused more people to search flu-related terms
- Is it time to recalibrate the model and/or hybridise both data sources?
- In 2014, Google retrained and updated the model significantly  
Will wait and see!

# Google flu trends, so what went wrong?

- Google had to mine 50 million search terms that best correlated with 1,152 “true” data points collected by CDC (Centers for Disease Control and Prevention). They retained 45 queries
- Winter is coming: Correlation is not causation
- Are these correlates stable and comparable over time?
- Google’s search algorithm changes very frequently. Autosuggest feature might have implicitly caused more people to search flu-related terms
- Is it time to recalibrate the model and/or hybridise both data sources?
- In 2014, Google retrained and updated the model significantly  
Will wait and see!

## Google flu trends, so what went wrong?

- Google had to mine 50 million search terms that best correlated with 1,152 “true” data points collected by CDC (Centers for Disease Control and Prevention). They retained 45 queries
- Winter is coming: Correlation is not causation
- Are these correlates stable and comparable over time?
- Google’s search algorithm changes very frequently. Autosuggest feature might have implicitly caused more people to search flu-related terms
- Is it time to recalibrate the model and/or hybridise both data sources?
- In 2014, Google retrained and updated the model significantly  
**Will wait and see!**

# Laws of data analysis

- ① Shite in, shite out  
*Anonymous*
- ② If you torture the data long enough it will confess to anything  
*Ronald Coase (1910 - 2013)*
- ③ A sufficiently elaborate analysis process can always lend an air of legitimacy  
*Chris Laws, ex-McLaren boss*

# Laws of data analysis

- ① Shite in, shite out  
*Anonymous*

- ② If you torture the data long enough it will confess to anything  
*Ronald Coase (1910 - 2013)*

- ③ A sufficiently elaborate analysis process can always lend an air of legitimacy  
*Chris Laws, ex-McLaren boss*



# Laws of data analysis

- ① Shite in, shite out  
*Anonymous*
- ② If you torture the data long enough it will confess to anything  
*Ronald Coase (1910 - 2013)*
- ③ A sufficiently elaborate analysis process can always lend an air of legitimacy  
*Chris Laws, ex-McLaren boss*

# Laws of data analysis

- ① Shite in, shite out  
*Anonymous*
- ② If you torture the data long enough it will confess to anything  
*Ronald Coase (1910 - 2013)*
- ③ A sufficiently elaborate analysis process can always lend an air of legitimacy  
*Chris Laws, ex-McLaren boss*

## A few important tips

- **Data exploration** can reveal important characteristics in the data e.g histograms, boxplots, scatter plot matrix (plot your raw data!)
- Covariates with **little variability** should be discarded
- How are the predictors **correlated** to each other?
- Missing data - ignore or **impute**?
- **Standardise** or **normalise** predictors with widely varying ranges
- **Features** extracted from data are **key**, they need to be directly relevant to the question you're asking
- Use **expert application knowledge** over automatic feature extraction
- **Never** trust predictions for inputs outside the training dataset range (i.e avoid extrapolation)

# Which machine learning algorithm should I use?

- A personal choice/what you're comfortable using rather than some rules set in stone
- Always start with simple models before using more complex ones
- Some methods are more appropriate than others in certain domains:

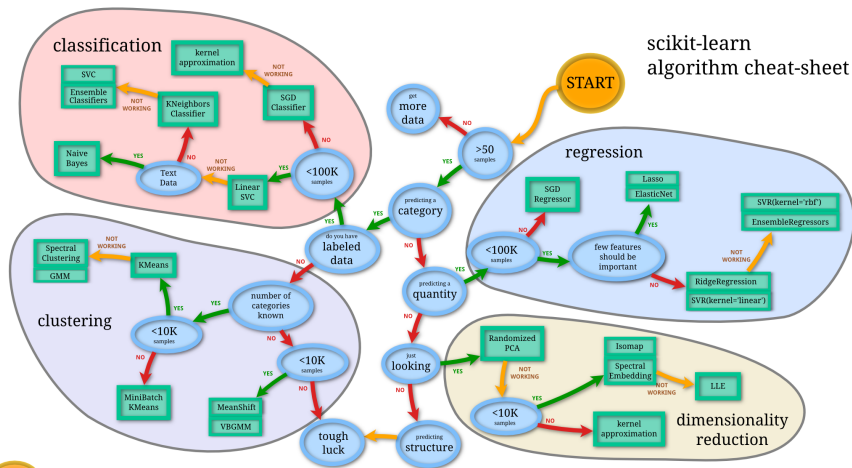
Interpretability: Decision trees or association rule learning

Lots of independent features and data: Naïve bayes

Knowledge of which features are correlated with each other: Bayesian network

Thousands of mixed categorical and continuous variables: Random forests

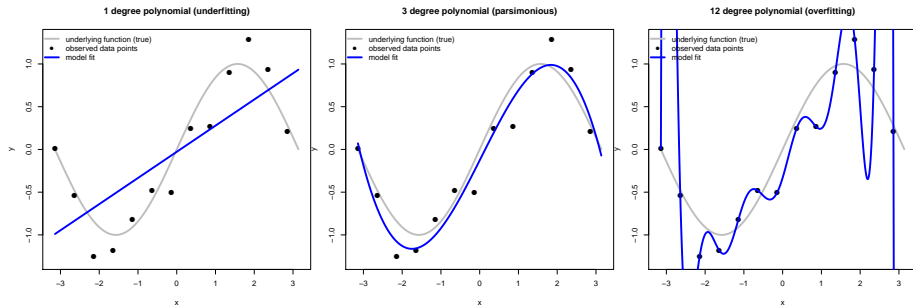
# Which machine learning algorithm should I use?



Source: [scikit-learn.org](http://scikit-learn.org)

# Selecting and validating models

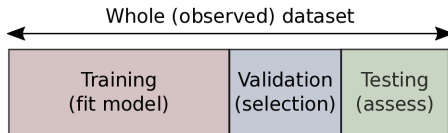
Recall the bias-variance tradeoff:



**How do we choose model complexity (model selection)?**  
**Once we choose it, how do we validate/assess it?**

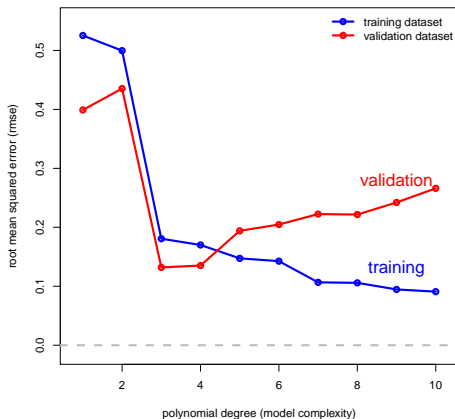
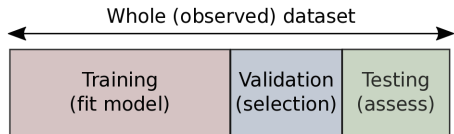
# Selecting and validating models

- Split data into three parts:



# Selecting and validating models

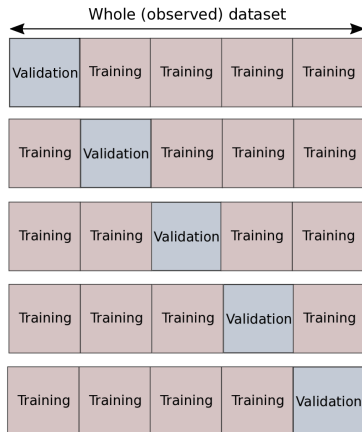
- Split data into three parts:





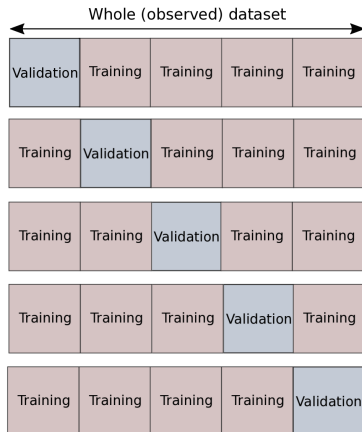
# Selecting and validating models

- When data is limited use  $k$ -fold cross-validation to estimate test error



# Selecting and validating models

- When data is limited use  $k$ -fold cross-validation to estimate test error



**Remember:** the most rigorous assessment is to apply the model to an “unseen” dataset